



Atty. Dkt. No. 043034/0158

0230

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

#3
12-28-00
9M

Applicant: Hidehito KUBO

Title: LOAD BALANCING METHOD AND SYSTEM BASED ON ESTIMATED
ELONGATION RATES

Appl. No.: 09/680,517

Filing Date: 10/06/2000

Examiner: Unassigned

Art Unit: Unassigned

RECEIVED

DEC 11 2000

Technology Center 2100

CLAIM FOR CONVENTION PRIORITY

Assistant Commissioner for Patents
Washington, D.C. 20231

Sir:

The benefit of the filing date of the following prior foreign application filed in the following foreign country is hereby requested, and the right of priority provided in 35 U.S.C. § 119 is hereby claimed.

In support of this claim, filed herewith is a certified copy of said original foreign application:

- Japanese Patent Application No. 11-285271 filed October 6, 1999.

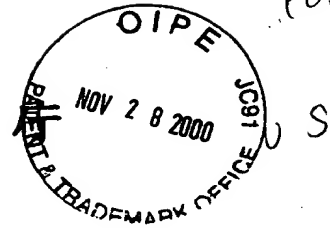
Respectfully submitted,

Date November 28, 2000

FOLEY & LARDNER
Washington Harbour
3000 K Street, N.W., Suite 500
Washington, D.C. 20007-5109
Telephone: (202) 672-5407
Facsimile: (202) 672-5399

By Philip J. Artavia Reg. No. 38,819
for / David A. Blumenthal
Attorney for Applicant
Registration No. 26,257

日 本 国 特 許
PATENT OFFICE
JAPANESE GOVERNMENT



別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

1999年10月 6日

出 願 番 号

Application Number:

平成11年特許願第285271号

出 願 人

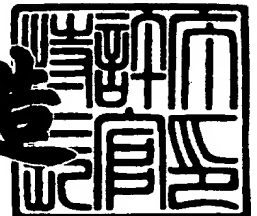
Applicant (s):

日本電気株式会社

2000年 7月14日

特許庁長官
Commissioner,
Patent Office

及川耕造



出証番号 出証特2000-3055214

【書類名】 特許願

【整理番号】 33509609

【提出日】 平成11年10月 6日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 15/16

【発明者】

 【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

 【氏名】 久保 秀士

【特許出願人】

 【識別番号】 000004237

 【氏名又は名称】 日本電気株式会社

【代理人】

 【識別番号】 100088959

 【弁理士】

 【氏名又は名称】 境 廣巳

【手数料の表示】

 【予納台帳番号】 009715

 【納付金額】 21,000円

【提出物件の目録】

 【物件名】 明細書 1

 【物件名】 図面 1

 【物件名】 要約書 1

 【包括委任状番号】 9002136

【ブルーフの要否】 要

【書類名】 明細書

【発明の名称】 推定伸長率に基づくトランザクション負荷分散方法及び方式並びにコンピュータ可読記録媒体

【特許請求の範囲】

【請求項 1】 トランザクション処理要求を発生する端末装置群と、該要求の処理を負荷分担して実行する複数の計算機からなるシステムにおいて、

各計算機の負荷状況を推定し、該推定負荷状況に基づいてすべての計算機について処理時間の推定伸長率を求め、この推定伸長率をベースとして各計算機の負荷指標の値を計算し、該負荷指標の値に基づいてトランザクション実行の各計算機への配分を決定することを特徴とする推定伸長率に基づくトランザクション負荷分散方法。

【請求項 2】 前記各計算機の負荷状況の推定に当たっては、一定時間ごとに各計算機の負荷データを測定してこれを元に各計算機の負荷状況を推定することを特徴とする請求項 1 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 3】 前記各計算機の負荷状況の推定に当たっては、各計算機におけるトランザクション処理の開始・終了に応じて各計算機の処理中トランザクション現在数を常に把握し、前記推定伸長率を求めるに当たっては、前記測定に基づいた推定負荷状況と前記処理中トランザクション現在数とに基づくことを特徴とする請求項 2 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 4】 端末からトランザクション処理要求が到着したときに、前記推定負荷状況に基づいてすべての計算機について前記処理時間の推定伸長率を求め、この推定伸長率をベースに各計算機の前記負荷指標の値を計算し、該到着した処理要求の実行に最適な計算機を該負荷指標の値から決定して処理させるようにすることを特徴とする請求項 1、2 または 3 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 5】 前記一定時間ごとに測定する各計算機の負荷データとして、処理中トランザクション数および CPU 系に滞在する業務処理プロセス数を測定して使用することを特徴とする請求項 2、3 または 4 に記載の推定伸長率に基づ

くトランザクション負荷分散方法。

【請求項 6】 前記一定時間ごとに測定する各計算機の負荷データとして、処理中トランザクション数および CPU 使用率を測定して使用することを特徴とする請求項 2、3 または 4 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 7】 前記各計算機の負荷状況の推定に当たって、一定時間ごとに測定した前記各計算機の負荷データの系列を総合的に用いて推定することを特徴とする請求項 2、3、4、5 または 6 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 8】 前記処理時間の推定伸長率を求めるにあたり、既に得られている前記推定負荷状況のデータを、前記処理中トランザクション現在数を用いて補正して使用することを特徴とする請求項 3、4、5、6 または 7 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 9】 前記計算機の負荷指標として、該計算機へ新規にトランザクションを割当てる前あるいは割当て後における、総推定伸長率、すなわち、該計算機における前記処理時間の推定伸長率に該計算機の処理中トランザクション数を乗じた値を用いることを特徴とする請求項 1、2、3、4、5、6、7 または 8 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 10】 前記計算機の負荷指標として、該計算機へ新規にトランザクションを割り当てた後における前記総推定伸長率と、割当て前における前記総推定伸長率との差を用いることを特徴とする請求項 1、2、3、4、5、6、7 または 8 に記載の推定伸長率に基づくトランザクション負荷分散方法。

【請求項 11】 トランザクション処理要求を発生する端末装置群と該要求の処理を負荷分担して実行する複数の計算機からなるシステムにおいて、

各計算機の負荷状況を推定する負荷データ測定手段と、

推定した該負荷状況を記憶する負荷データ記憶手段と、

該推定負荷状況に基づいてすべての計算機について処理時間の推定伸長率を求め、この推定伸長率をベースとして各計算機の負荷指標の値を計算し、該負荷指標の値に基づいてトランザクション実行の各計算機への配分を決定する実行計算

機選択手段と、

前記各計算機ごとにその上に存在し、複数のトランザクション実行を並列に行い、前記実行計算機選択手段に指令されたトランザクションの実行を管理するトランザクション処理手段とを備えたことを特徴とする推定伸長率に基づくトランザクション負荷分散方式。

【請求項 1 2】 前記負荷データ測定手段は、一定時間ごとに各計算機の負荷データを測定してこれを元に各計算機の負荷状況を推定することを特徴とする請求項 1 1 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 1 3】 前記負荷データ測定手段は、各計算機におけるトランザクション処理の開始・終了に応じて各計算機の処理中トランザクション現在数を常に把握し、

前記実行計算機選択手段は、すべての計算機について前記処理時間の推定伸長率を求めるに際し、前記測定に基づいた推定負荷状況と前記処理中トランザクション現在数とに基づくことを特徴とする請求項 1 2 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 1 4】 前記実行計算機選択手段は、端末から前記トランザクション処理要求が到着したときに起動され、該到着した処理要求の実行に最適な計算機を決定して処理させるようにすることを特徴とする請求項 1 1、1 2 または 1 3 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 1 5】 前記負荷データ測定手段は、前記一定時間ごとに測定する各計算機の負荷データとして、処理中トランザクション数および CPU 系に滞在する業務処理プロセス数を測定し、負荷状況の推定に使用することを特徴とする請求項 1 2、1 3 または 1 4 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 1 6】 前記負荷データ測定手段は、前記一定時間ごとに測定する各計算機の負荷データとして、処理中トランザクション数および CPU 使用率を測定し、負荷状況の推定に使用することを特徴とする請求項 1 2、1 3 または 1 4 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 1 7】 前記負荷データ測定手段は、前記各計算機の負荷状況の推

定に当たって、一定時間ごとに測定した前記各計算機の負荷データの系列を総合的に用いて推定することを特徴とする請求項 12、13、14、15 または 16 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 18】 前記実行計算機選択手段は、各計算機の前記処理時間の推定伸長率を求めるにあたり、既に得られている前記推定負荷状況のデータを前記処理中トランザクション現在数を用いて補正して使用することを特徴とする請求項 13、14、15、16 または 17 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 19】 前記実行計算機選択手段は、各計算機の前記負荷指標として、該計算機へ新規にトランザクションを割り当てる前あるいは割り当て後における、総推定伸長率、すなわち、該計算機における処理時間の推定伸長率に該計算機の処理中トランザクション数を乗じた値、を用いることを特徴とする請求項 11、12、13、14、15、16、17 または 18 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 20】 前記実行計算機選択手段は、各計算機の前記負荷指標として、該計算機へ新規にトランザクションを割り当てた後における前記総推定伸長率と、割り当て前における前記総推定伸長率との差を用いることを特徴とする請求項 11、12、13、14、15、16、17 または 18 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 21】 前記実行計算機選択手段と前記負荷データ記憶手段とがそれぞれシステムに一つだけ存在して集中的にその機能を実行し、前記実行計算機選択手段は各計算機の前記負荷指標の値を直接的に反映してトランザクションの配分を行うことを特徴とする請求項 11、12、13、14、15、16、17、18、19 または 20 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 22】 前記実行計算機選択手段が各計算機ごとに一つずつ分散して存在し、固定的、静的または準静的な分配方式によって計算機に配分されてきたトランザクションについて、該計算機上の実行計算機選択手段が、すべての計算機の前記負荷指標の値に基づき、次の 2 つの決定すなわち、該計算機でそのま

ま処理するか他に回すかを閾値判断で決定、および他に回す場合はその送付先を決定、を行うことを特徴とする請求項 1 1、1 2、1 3、1 4、1 5、1 6、1 7、1 8、1 9 または 2 0 に記載の推定伸長率に基づくトランザクション負荷分散方式。

【請求項 2 3】 トランザクション処理要求を発生する端末装置群と、該要求の処理を負荷分担して実行する複数の計算機からなるシステムにおいて、各計算機の処理時間の伸長率を推定し、この推定伸長率をベースとした各計算機の負荷指標に基づいてトランザクション処理要求を各計算機へ配分することを特徴とする推定伸長率に基づくトランザクション負荷分散方法。

【請求項 2 4】 トランザクション処理要求を発生する端末装置群および該要求の処理を負荷分担して実行する複数の計算機に接続され、前記端末装置群からのすべての処理要求を集中的に受け取って前記計算機に配分する中継配分装置を構成するコンピュータに、

前記各計算機から一定時間ごとに通知される当該計算機の CPU 系に滞在する業務処理プロセス数あるいは CPU 使用率、および処理中トランザクション数を含む負荷データに基づき、各計算機の負荷状況を推定するステップ、

推定された負荷状況に基づいてすべての計算機について処理時間の推定伸長率を求めるステップ、

該推定伸長率をベースとして各計算機の負荷指標の値を計算するステップ、

該負荷指標の値に基づいてトランザクション実行の各計算機への配分を決定するステップ、

を実行させるプログラムを記録したコンピュータ可読記録媒体。

【請求項 2 5】 トランザクション処理要求を発生する端末装置群からの要求の処理を負荷分担して実行する複数の計算機のそれぞれに、

自計算機の CPU 系に滞在する業務処理プロセス数あるいは CPU 使用率、および処理中トランザクション数を含む負荷データを一定時間ごとに測定し、これを元に自計算機の負荷状況を推定して記憶すると共に他のすべての計算機に通知するステップ、

処理要求の到着を契機に、すべての計算機についての前記推定負荷状況から各

計算機の処理時間の推定伸長率、該推定伸長率をベースにした各計算機の負荷指標を求めるステップ、

該求めた各計算機の負荷指標の値に基づいて前記到着した処理要求を自計算機で処理すべきか他計算機に依頼すべきかを判断し、他に送付する場合はその送り先を決定し、選択した実行先に送付依頼するステップ、
を実行させるプログラムを記録したコンピュータ可読記録媒体。

【請求項 2 6】 トランザクション処理要求を発生する端末装置群からの処理要求を一括して受け取って静的／準静的な方式により当該処理要求を配分する計算機を決定する中継仮配分装置に接続された前記それぞれの計算機に、

自計算機の CPU 系に滞在する業務処理プロセス数あるいは CPU 使用率、および処理中トランザクション数を含む負荷データを一定時間ごとに測定し、これを元に自計算機の負荷状況を推定して記憶すると共に他のすべての計算機に通知するステップ、

処理要求の到着を契機に、すべての計算機についての前記推定負荷状況から各計算機の処理時間の推定伸長率、該推定伸長率をベースにした各計算機の負荷指標を求めるステップ、

該求めた各計算機の負荷指標の値に基づいて前記到着した処理要求を自計算機で処理すべきか他計算機に依頼すべきかを判断し、他に送付する場合はその送り先を決定し、選択した実行先に送付依頼するステップ、
を実行させるプログラムを記録したコンピュータ可読記録媒体。

【発明の詳細な説明】

【0 0 0 1】

【発明の属する技術分野】

本発明は、クラスタ構成などと呼ばれる比較的緊密に結合された複数の計算機が負荷分担してトランザクション処理を実行するシステムにおいて、トランザクション処理の負荷を各計算機に分散させる方式に関し、特に各計算機における負荷状況を示す指標である「処理時間の推定伸長率」に基づいて処理要求を動的に配分することにより計算機間の負荷をバランスさせ、全体として応答時間の平均及びばらつきを小さく保つ方式に関する。

【0002】

【従来の技術】

この種の負荷分散方式は、複数の処理装置（計算機）をもつシステムにおいて、規模の小さい処理を要求するメッセージが次々に大量に到着するのを、これら複数の計算機に適切に分配することによって計算機間で負荷を分散させ、システムから最大の性能を引き出そうとする。個々の処理が小規模なので、一般に、処理の途中で移動させることは考えず、到着時に処理を実行すべき計算機を決定してそこへ送付し、終了までそこで処理させる。また、対話型処理なので、負荷分散の最終的な目標は応答時間の平均（及びばらつき）の最小化である。処理要求の到着時に実行させる計算機を決定するが、この実行計算機を選択に当たっては、基本的に、負荷が最も低い計算機を選ぶことになる。ここで、何を「負荷の指標」とするかという問題が生ずる。従来、負荷の指標としては、計算機のCPU使用率、実行中の処理数、近い過去の応答時間の実績などが、個々に単独で、あるいは組み合わせて用いられていた。

【0003】

従来のシステムの一例が、特開平10-312365号公報に記載されている（従来技術1と呼ぶ）。ここでは、一定時間ごとにサーバー（計算機）の負荷状態（実施例によればCPU使用率）を計測して格納しておき、処理要求到着時には格納してある負荷状態に基づき最も負荷の低いサーバーを実行計算機として決定する。また、端末側で応答時間を監視していて、これが所定値を越えているなら、経路変換（実行計算機の変更）要求を出す。ここでは、第一の負荷指標としてCPU使用率が用いられている。CPU使用率はよい指標ではあるが、測定値は過去の一定時間間隔における平均値であり、その後にあった処理の開始終了などの影響を含まないこともあり、「現在の負荷」を表しているという意味での信頼性はあまり高くない。特に、動的な制御の下では、次の測定時までは同じデータが用いられるので、この間に到着した処理要求は負荷が最小であると判断された特定の一つのサーバーに集中的に送られることになり、負荷のシーソー現象を起こす可能性がある。第二の負荷指標として、実行途中の処理の応答時間の実績が、その処理自身の実行先切り替えの判断に用いられている。ここにおけるよう

に実行途中における実行計算機の切り替えがオーバーヘッド小さく可能な場合があるとすれば、有効な指標であろう。しかし、到着時の配分においては、その処理自身の実績がないのでこの指標は使用不可能である。

【0004】

従来のシステムの別の一例が、特開平10-27168号公報に記載されている（従来技術2と呼ぶ）。ここでは、各計算機で最後に実行を終了したメッセージについてその処理時間を記憶しておき、この時間にその計算機上で処理中のメッセージの数を乗じたものを負荷の指標として用いている。メッセージの到着時には、すべての計算機についてこの負荷指標の値を計算し、この値の最も小さい計算機にメッセージを送り処理を依頼する。この場合、最後に終了したメッセージの処理時間がその計算機上の処理時間を代表するか、という問題がある。この処理時間は、その計算機の混み具合とその最終終了メッセージ処理のジョブ特性（純処理時間、CPU／入出力の比率）とを反映しているはずである。すべてのメッセージについて後者のジョブ特性が同一であるならば、負荷指標として目安を与えられと考えられる。しかし、様々なジョブ特性のものが混在している一般の状況では、個々の処理時間実績をそのまま負荷状態を反映するものと考えたと判断を誤る可能性が大きい。

【0005】

他の従来のシステムの例が、特開平7-302242号公報に記載されている（従来技術3と呼ぶ）。ここには多くの請求項があるが、本発明に近いものは請求項9、段落番号「152」～「161」に記載されたものである。ここでは、トランザクション処理部の負荷を定期的に検出して時刻と共に負荷の履歴を記憶しておく。そして、負荷傾向 T_r を次の式で計算する。

$$T_r = (W_2 - W_1) / (T_2 - T_1)$$

トランザクション処理要求を受けると、一定時間 T_i 後に処理負荷予測値が閾値 W_t を越えないと判断した場合（ $T_r \cdot T_i \leq W_t$ ）、自分で受け付け、そうでなければ拒否する。あるいは、より負荷の低い他のサーバに処理を依頼する。この例では、負荷を定期的に検出してこれをベースとして判断するが、具体的に何を「負荷」とするかについては公報全体を通じて明確に規定されていない

。負荷の指標を規定することは負荷配分にとって重要な第一歩であるが、それがなされていない。また、 T_i 後の負荷を過去の線形外挿により予測しようとしているが($T_r \cdot T_i$ では不足と思われるが、それは別として)、これは良い予測とは思えない。システム全体の負荷についてはこの種のマクロな予測も有効かも知れないが、自サーバについてはその先の負荷状況は、現在の状態と処理中トランザクションの終了タイミング、自分が新たに処理を受け入れるか否かというミクロな動きで決まるものであり、過去の傾向を延長してそのまま信じてしまうのは危険である。

【0006】

また、上述の従来技術1, 2, 3はいずれも、到着したメッセージ自体の処理時間を最短にすることを狙って実行先を決定している。しかし、このような個別最適化がシステム全体としての最適化に直結するという保証は必ずしもあるわけではない。

【0007】

他の従来のシステムの例が、1997年にSpringer社から発行されたOptimal Load Balancing in Distributed Computer Systems (H.Kameda他著) の第225頁～第232頁に記載されている(従来技術4と呼ぶ)。前提としているモデルを1台のCPUに注目して示すと図2のようなものである。ジョブは到着すると、CPU(図では計算機*i*)とディスクの使用を繰り返し、処理を終了すると立ち去る。この間の時間が応答時間である。複数のジョブが並行処理されるので、CPUの前には待ち行列が生ずる。他のCPUも、図示した計算機*i*と同じ位置づけになり、ディスクに対するアクセス時間はすべてのCPUから同等である。ここでは、負荷指標として次の2つの式が示されている。

$$f_i = s_i (n_i + 1)^2 \quad (\text{式1})$$

$$F_i = s_i (n_i + 1) \quad (\text{式2})$$

ここで、 f および F は負荷指標、 i は計算機番号、 s はジョブのCPUにおける純サービス時間の平均、 n はCPU系に存在するジョブ数である。これらの式は、待ち行列理論で言う開放型待ち行列網モデルにおいて、CPU系について平衡状態における平均値に関して成立する関係をもとに、平均応答時間がある意味で

最小化するという目的のために、小さいほどよい値として導き出されたものである。実際、(式2)はCPU系における平均滞在時間を表し、これに入出力系における平均滞在時間を加えると平均応答時間となるものである。これらの負荷指標は、静的負荷配分のための指標としてはある意味の最適性が証明されている。しかし、動的制御はその時々状況に応じた制御を行いうるところにメリットがあり、 s_i 、 n_i については、平衡状態における平均値でなく現在値を用いないと意味がない。現在値に関し n_i は測定可能であるが、計算機*i*上で実行中のジョブミックスの特性を反映する s_i は直接には測定不可能である。当文献上における評価では、 s_i として全体の平均値を用いている。実行する処理が、ジョブ特性の観点から一種類でしかも特性のばらつきが小さいなら全体の平均値を用いてもよいであろうが、一種類でもばらつきが大きい場合や、現実には一般的と考えられる特性の異なる複数種類の処理が混在する場合には、全体の平均値を用いてしまうと動的制御のメリットが大きく失われることになる。

【0008】

【発明が解決しようとする課題】

第1の問題点は、負荷分散のベースとなる各計算機の負荷状況の把握が不十分であるということである。従来は、処理中トランザクション数、CPU使用率などがそれぞれ単独であるいは組み合わせて負荷指標として使われていたが、これらは、その時点で処理中のトランザクション群のCPU/入出力使用比率を含むジョブ特性まで含めた、システムの混み具合を十分に反映するものとは言えない。また、トランザクション処理のように小規模の処理要求が大量に到着するシステムでは、短い時間間隔で正確に負荷状況を把握する必要があるが、分散システムにおけるデータ収集のオーバーヘッドを恐れ、収集頻度を少なくするような傾向があった。低オーバーヘッドな良質のデータを用いる工夫と共に、クラスタ型などの環境ではデータ収集は高速・低オーバーヘッドなので、これを生かして良い負荷分散を実現するような方式が求められる。

【0009】

第2の問題点は、必ずしも、システム全体としての最適化(応答時間の平均、分散の最小化)を図るものではなかったということである。到着したトランザク

ションの配分先を決定するに際し、当該トランザクションにとってその時点で最適な（最短時間で処理できると予想される）計算機を選択していたが、このような個別最適化は、システム全体としての最適化につながることを、必ずしも保証するものではない。

【0010】

【発明の目的】

本発明の目的は、クラスタ構成などと呼ばれる比較的緊密に結合された複数の計算機が負荷分担してトランザクション処理を実行するシステムにおいて、到着した処理要求に対し適切な実行先計算機を選択するための基準となる有効な負荷指標を提供し、これに基づく選択を小さいオーバーヘッドで実行可能にすることにより、トランザクション処理の負荷を短期レンジでも計算機間でバランスさせ、もって、全体として応答時間の平均とばらつきを小さく保つことを可能にする動的な負荷分散方式を提供することにある。

【0011】

【課題を解決するための手段】

本発明は、トランザクション処理要求を発生する端末装置群と該要求の処理を負荷分担して実行する複数の計算機からなるシステムにおいて、各計算機の処理時間の伸長率を推定し、この推定伸長率をベースとした各計算機の負荷指標に基づいてトランザクション処理要求を各計算機へ配分する。具体的には、本発明にかかる負荷分散方法にあつては、各計算機の負荷状況を推定し、該推定負荷状況に基づいてすべての計算機について処理時間の推定伸長率を求め、この推定伸長率をベースとして各計算機の負荷指標の値を計算し、該負荷指標の値に基づいてトランザクション実行の各計算機への配分を決定する。また、本発明にかかる負荷分散方式にあつては、各計算機の負荷状況を推定する負荷データ測定手段と、推定した該負荷状況を記憶する負荷データ記憶手段と、該推定負荷状況に基づいてすべての計算機について処理時間の推定伸長率を求め、この推定伸長率をベースとして各計算機の負荷指標の値を計算し、該負荷指標の値に基づいてトランザクション実行の各計算機への配分を決定する実行計算機選択手段と、前記各計算機ごとにその上に存在し、複数のトランザクション実行を並列に行い、前記実行計

算機選択手段に指令されたトランザクションの実行を管理するトランザクション処理手段とを有する。

【0012】

各計算機の処理時間の伸長率とは、業務処理プロセスの応答時間、すなわち待ち時間も含む処理時間の、純処理時間（CPU、ファイル装置という資源を実際に使用する時間の合計）に対する倍率を意味する。この伸長率は、当該計算機上で実行中のプロセスの集まり（ジョブミックス）の、動作中の群としてのプログラム特性（CPU使用特性だけでなく、CPU-I/O使用特性を含む）を反映している。従って、処理速度が同じ計算機ならば、同一の処理は伸長率の小さい計算機で実行した方が処理時間は短くなり、応答時間を短くできる。

【0013】

各計算機における処理時間の伸長率は、一定時間ごとに各計算機の負荷データとして例えば処理中トランザクション数とCPU系に滞在する業務処理プロセス数、または、処理中トランザクション数とCPU使用率を測定し、これらに基づいて推定する。一定時間ごとに測定した負荷データの系列を総合的に用いて各計算機の負荷状況を推定したり、各計算機におけるトランザクション処理の開始・終了に応じて各計算機の処理中トランザクション現在数を常に把握しておき、この処理中トランザクション現在数を用いて推定負荷状況データを補正したりすれば、推定負荷状況の推定精度が高まり、ひいては伸長率の推定精度も高まる。

【0014】

推定伸長率をベースとして各計算機の負荷指標の値を求め、これに応じて、到着する処理要求を計算機へスケジュールする。負荷指標としては、推定伸長率そのものを負荷指標とすることができる他、当該計算機へ新規にトランザクションを割当て前あるいは割当て後における、総推定伸長率、すなわち、当該計算機における前記処理時間の推定伸長率に当該計算機の処理中トランザクション数を乗じた値を用いることができ、また、当該計算機へ新規にトランザクションを割り当てた後における前記総推定伸長率と、割当て前における前記総推定伸長率との差を用いることもできる。後者では、伸長率の増分最小の計算機が選択されるため、システム全体にとって当スケジュールによる応答時間総和の増加を最小に

する選択になる。

【 0 0 1 5 】

本発明の推定伸長率に基づくトランザクション負荷分散方式では、実行計算機選択手段と負荷データ記憶手段とがそれぞれシステムに一つだけ存在して集中的にその機能を実行し、前記実行計算機選択手段は各計算機の前記負荷指標の値を直接的に反映してトランザクションの配分を行うよう構成して良い。具体的には、すべての処理要求を集中的に受け取って計算機に配分する中継配分装置（図 1 の 2）を備え、各計算機（図 1 の 1 x）上に存在して一定時間ごとに負荷データを測定し中継配分装置に通知する手段（図 1 の 1 x 1 と 1 x 3）と、これを受けて中継配分装置上で各計算機の負荷状況を推定して記憶すると共に各計算機の処理中トランザクション現在数を常に把握する手段（図 1 の 8 と 6）と、中継配分装置上に存在し端末から処理要求が到着すると起動され、該到着処理要求を処理する計算機を決定して送付する実行計算機選択手段（図 1 の 7）とを備え、実測に基づいたその時点の推定負荷状況と処理中トランザクション現在数とに基づいて各計算機の処理時間の推定伸長率を求め、この推定伸長率をベースに各計算機の負荷指標の値を計算し、該負荷指標の値から、到着した処理要求に関して動的負荷配分の観点から最適な計算機を決定して処理させるように動作する（この方式を第 1 の方式と呼ぶ）。

【 0 0 1 6 】

本発明の推定伸長率に基づくトランザクション負荷分散方式では、また、実行計算機選択手段が各計算機ごとに一つずつ分散して存在し、固定的、静的または準静的な分配方式によって計算機に配分されてきたトランザクションについて、該計算機上の実行計算機選択手段が、すべての計算機の前記負荷指標の値に基づき、次の 2 つの決定すなわち、該計算機でそのまま処理するか他に回すかを閾値判断で決定、および他に回す場合はその送付先を決定、を行うよう構成して良い。具体的な構成例としては、次の 2 つの方式が考えられる（それぞれ第 2 の方式、第 3 の方式と呼ぶ）。

【 0 0 1 7 】

第 2 の方式は、中継配分装置を備えず、したがって処理要求は静的な方式で各

計算機へ配分されるが、各計算機上に存在して一定時間ごとに負荷データを測定しこれを元に負荷状況を推定して記憶すると共に他のすべての計算機に通知する手段（図 6 の 1 x 1 と 8 x）と、処理要求の到着と共に到着計算機上で起動されてすべての計算機についての前記推定負荷状況から各計算機の処理時間の推定伸長率を求めこれを元に負荷指標の値を計算し、各計算機の該負荷指標の値に基づいて自計算機で処理すべきか他計算機に依頼すべきかを閾値を用いて判断し、他に送付する場合はその送り先を決定し、選択した実行先に送付依頼するよう動作する実行計算機選択手段（図 6 の 7 x）とを備える構成である。

【 0 0 1 8 】

第 3 の方式は、処理要求を一括して受け取って静的／準静的方式により計算機に配分する中継仮配分装置（図 7 の 2 5）を備え、動的に最適とは言えない仮配分がなされるが、各計算機上に第 2 の方式と同一の、一定時間ごとに働く測定手段と、負荷状況推定手段、及び仮配分された処理要求の到着時に起動され各計算機の推定伸長率を計算しこれに基づいて負荷指標の値を求め、自計算機で処理すべきか他に依頼するとしたらどの計算機かを決定し処理を依頼するように動作する実行計算機選択手段（図 7 の 7 x）を備える構成である。

【 0 0 1 9 】

【発明の実施の形態】

次に、本発明の実施の形態について図面を参照して詳細に説明する。

【 0 0 2 0 】

図 1 を参照すると、本発明の第 1 の実施の形態は、プログラム制御により動作する計算機群 1 と、中継配分装置 2 と、高速チャネル 3 と、ファイル装置群 4 と、端末装置群 5 と、通信網 5 1 とから構成されている。

【 0 0 2 1 】

計算機群 1 は計算機 1 1 ～ 1 n を含み、計算機 1 1 ～ 1 n は、それぞれ、負荷データ測定手段 A 1 1 1 ～ 1 n 1 と、トランザクション処理手段 1 1 2 ～ 1 n 2 と、通信手段 1 1 3 ～ 1 n 3 とを含み、トランザクション処理手段 1 1 2 ～ 1 n 2 はそれぞれ複数の業務処理プロセスを含み（図示せず）、中継配分装置 2 は通信手段 2 1 と、負荷データ記憶手段 6 と、実行計算機選択手段 7 と、負荷データ

測定手段 B 8 とを含む。計算機 1 1 ~ 1 n は主記憶を共有しないが、ファイル装置群 4 には高速チャンネル 3 を介して性能的に同等の条件で接続されており、ファイル装置を共有している。中継配分装置 2 も高速チャンネル 3 を介して計算機群 1 に接続され、端末群 5 は通信網 5 1 を介して中継配分装置 2 に接続されている。

【 0 0 2 2 】

1 つのトランザクションの処理の概略は次のようになる。トランザクション処理要求であるメッセージは端末群 5 に属する端末装置から送り出され、中継配分装置 2 に伝えられる。中継配分装置 2 は、受け取ったメッセージを処理する計算機 1 i を決定し、該計算機に高速チャンネル 3 を介して該メッセージを送る。受け取った計算機 1 i ではトランザクション処理を実行し、応答メッセージを作成して逆の経路を通して要求元端末に返す。

【 0 0 2 3 】

ここで、各計算機 1 1 ~ 1 n に備わる負荷データ測定手段 A 1 1 1 ~ 1 n 1、トランザクション処理手段 1 1 2 ~ 1 n 2、通信手段 1 1 3 ~ 1 n 3 をソフトウェア的に実現する場合、これらの各手段を実現するプログラムは図示しない C D - R O M、磁気ディスク、半導体メモリ等の機械読み取り可能な記録媒体に保存されており、計算機群 1 の立ち上げ時などに記録媒体に記録されたプログラムが各計算機に読み込まれ、各計算機の動作を制御することにより、各計算機上にこれら各手段を実現する。また、中継配分装置 2 に備わる通信手段 2 1、負荷データ記憶手段 6、実行計算機選択手段 7、負荷データ測定手段 B 8 をソフトウェア的に実現する場合、これらの各手段を実現するプログラムは図示しない C D - R O M、磁気ディスク、半導体メモリ等の機械読み取り可能な記録媒体に保存されており、中継配分装置 2 を構成する計算機の立ち上げ時などに記録媒体に記録されたプログラムがその計算機に読み込まれ、その計算機の動作を制御することにより、その計算機上にこれら各手段を実現する。

【 0 0 2 4 】

上記の各手段はそれぞれ概略つぎのように動作する。

【 0 0 2 5 】

計算機 1 1 ~ 1 n の各々は同一の機能を持つので、以下では 1 i で代表させる

。計算機 1 i 上の負荷データ測定手段 A 1 i 1 は、一定時間ごとに自身の属する計算機 1 i の負荷データを測定し、結果を中継配分装置 2 に送る。トランザクション処理手段 1 i 2 は、中継配分装置 2 から送られた処理要求メッセージを通信手段 1 i 3 から受け取ると、自身の管理下の業務処理プロセスを割り当て、処理を行わせる。前記業務処理プロセスは、プログラムの実行のために CPU の使用とファイル装置 4 上のファイルへのアクセスを繰り返し、処理が終了すると応答メッセージを作成し、通信手段 1 i 3 を介して中継配分装置 2 に送る。

【 0 0 2 6 】

中継配分装置 2 上の負荷データ記憶手段 6 には各計算機の負荷データが格納されている。負荷データ測定手段 B 8 は、各計算機から一定時間ごとに送られてくる負荷データを通信手段 2 1 経由で受け、これを加工して推定データとして前記負荷データ記憶手段 6 に格納する。また、計算機へのトランザクション処理要求送付および応答メッセージ到着の通知を通信手段 2 1 から受け、負荷データ記憶手段 6 上の一部のデータを更新する。実行計算機選択手段 7 は、前記処理要求メッセージを通信手段 2 1 から渡され、負荷データ記憶手段 6 に記憶されている各計算機の前記推定負荷データから各計算機における推定伸長率を求め、これに基づいて実行すべき計算機を決定し、該計算機に向けて、通信手段 2 1 に処理要求メッセージを送付させる。

【 0 0 2 7 】

次に、図 1 ～図 5 を参照して本実施の形態の全体の動作について詳細に説明する。

【 0 0 2 8 】

図 1 において、トランザクション処理手段 1 i 2 は、中継配分装置 2 から送られた処理要求メッセージを通信手段 1 i 3 から受け取ると、自身の管理下の業務処理プロセスを割り当て、要求に応じたトランザクション処理を行わせる。トランザクション処理手段 1 i 2 は、複数のトランザクション、したがって複数の業務処理プロセスをマルチプログラミング状態で走らせることができ、これによって応答時間、資源使用効率を向上させている。業務処理プロセスは、適用業務プログラム実行のために CPU の使用とファイルアクセスのためのファイル装置群

4 への入出力を繰り返し、処理が終了すると応答メッセージを作成し、通信手段 1 i 3 を介して中継配分装置 2 に送りプロセスを終了する。

【0029】

図 2 に、1 台の計算機 i について、性能面から見たシステムのモデルを示す。業務処理プロセスの資源使用特性は CPU 使用時間と入出力の回数で捉えられるが、これはトランザクションごとに異なるものである。複数の処理を並行して走らせるので、資源の競合が起こる。そのため、少なくとも CPU の前にはプロセス待ち行列ができることが想定される。一般に、使用率の高い CPU ほど待ち時間が長い。入出力に関しては、いずれの計算機からも性能的に同条件にあるので、ここの処理時間は待ち時間も含めてアクセス元の計算機による差はないと考える。 N_i は計算機 i 上の処理中業務処理プロセス数を表すものとする。これは該計算機上で処理中のトランザクション数に相当する。 P_i は CPU 系に存在する業務処理プロセス数を表すものとする。これは N_i のうち CPU 割当て待ち（レディ状態）あるいは CPU 使用中であるプロセスの総数である。ファイル装置は共有されており、すべての計算機についてアクセス性能は同等なので、他の計算機も、性能的にはすべて計算機 i と同じ位置づけとなる。

【0030】

計算機 1 i 上の負荷データ測定手段 A 1 i 1 は、一定時間ごとに自身の属する計算機 1 i の負荷データとして、その時点で CPU 系に存在する前記業務処理プロセス数 P_i あるいは直前の測定以後今回までの間の CPU 使用率 R_i 、およびその時点での処理中業務処理プロセス数 N_i を測定し、結果を通信手段 1 i 3 を介して中継配分装置 2 に送る。前記測定の間隔は負荷分散の精度を左右するので、オーバーヘッドとの兼ね合いもあるが、通常のトランザクション処理では 100 ミリ秒程度以下、できれば 10 ミリ秒程度、であることが望ましい。 P_i を用いるか R_i を用いるかは実施システムごとに決定してよい。 P_i を用いる場合を P 方式、 R_i を用いる場合を R 方式と呼ぶことにする。

【0031】

図 3 に、中継配分装置 2 上の負荷データ記憶手段 6 に記憶する負荷データをテーブル形式で示す。計算機番号 T 1 はシステム内で稼働中の計算機の識別を示し

、テーブル上のデータは計算機ごとに1行を用意して管理されている。負荷データ測定手段B 8は、前記負荷データ測定手段A 1 i 1から一定時間ごとに送られるデータである、CPU系に存在する業務処理プロセス数P iあるいは直前の測定以後今回までの間のCPU使用率R i、およびその時点での処理中業務処理プロセス数N iの値を通信手段2 1経由で受ける。そして、これを加工して推定データとして、前記P iの推定値P e iあるいは前記R iの推定値R e i、および前記N iの推定値N e iを求め、前記R方式ならR e iからP e iを計算し、列T 3にN e iを列T 4にP e iをそれぞれi番目の値として格納する。

【0032】

測定値をそのまま用いず推定値に変換するのは、過去のデータを総合的に組み込むことによりサンプリングの信頼性の低さを補うためである。具体的な求め方として、次の方法がある。測定値をm、推定値をe、最新の測定がn回目であったとする。

$$e(n) = a * m(n) + (1 - a) * e(n - 1) \quad (\text{式3})$$

ここで、aはパラメタ ($0 < a \leq 1$) であり、また、e (n) の初期値e (0) はe (1) と等しいとする。すなわち、今回の測定値にaを乗じたものと前回の推定値に1 - aを乗じたものとの和を今回の推定値とする。式3は次のように展開できる。

$$e(n) = a * m(n) + a(1 - a) * m(n - 1) + a(1 - a)^2 * m(n - 2) + a(1 - a)^3 * m(n - 3) + \dots$$

この式は、推定値が、近い過去の測定値ほど重視する形で過去の測定値を全部取り込んだものになっていることを示している。aが大きい(1に近い)ほど近い過去を重視する度合いが高いことになる。前記測定間隔が十分に小さいなら、aの値は0.1などの小さい値とした方が推定値の信頼度は上がる。前記R方式では、式3により求められたR e iから次の式によりP e iを求め、列T 4に格納する。

$$P e i = R e i / (1.0 - R e i) \quad (R e i \geq 0.99 \text{ なら } P e i = N e i)$$

この式は、M/M/1待ち行列における系の長さ和使用率との関係そのものである。

【0033】

また、計算機 i 上で処理中のトランザクション現在数 N_{pi} が、テーブル上の列 $T2$ に保持されている。この値は、前記負荷データ測定手段 $B8$ が、計算機へのトランザクション処理要求送付（トランザクション開始）および応答メッセージ到着（トランザクション終了）の通知を通信手段 21 から受けて更新し、保持する。したがって、信頼できる測定値である。

【0034】

図 $4A$ は、本発明の第 1 の実施の形態の実行計算機選択手段 7 の動作を示すフローチャートである。実行計算機選択手段 7 は、端末から処理要求メッセージが中継配分装置 2 に到着すると通信手段 21 に起動されて該メッセージを渡され（図 $4A$ の 71 ）、負荷データ記憶手段 6 に記憶されている各計算機の負荷データから各計算機における推定伸長率を求めて（図 $4A$ の 72 、 73 ）負荷指標の値を計算し（図 $4A$ の 74 ）、これに基づいて実行すべき計算機を決定し（図 $4A$ の 75 ）、該計算機に向けて、通信手段 21 に処理要求メッセージを送付させる（図 $4A$ の 76 ）。

【0035】

推定値の補正（図 $4A$ の 72 ）では、負荷データ記憶手段 6 上のデータをベースに補正を行い、現時点における負荷データとして、処理中業務処理プロセス数（補正值） N_{ri} 、および CPU 系に存在する業務処理プロセス数（補正值） P_{ri} を次の式により求める。

$$N_{ri} = w * N_{pi} + (1 - w) * N_{ei} \quad (\text{式 } 4)$$

ここで w は重み係数（ $w \leq 1.0$ ）であり、 0.8 程度がよい。

$$\begin{aligned} P_{ri} &= P_{ei} + (N_{ri} - N_{ei}), \quad N_{ri} \geq N_{ei} \text{ のとき} \quad (\text{式 } 5) \\ &= P_{ei} * (N_{ri} / N_{ei}), \quad N_{ri} < N_{ei} \text{ のとき} \end{aligned}$$

式 5 は、過去のサンプリングから推定していた P_{ei} を、これと同一条件で推定した N_{ei} と現時点の状況を表す最も信頼できる補正值である N_{ri} との関係から、補正するものである。一定時間間隔でしか行わない推定をベースに、最新の推定値から現時点の真の値に近い補正值を得ることができる。

【0036】

推定伸長率の計算（図 4 A の 73）では、推定伸長率 E_{pi} を次の式により求め、結果をテーブルの列 T5 に格納する。

$X = N_{ri} * (P_{ri} + 1)$ として、

$$E_{pi} = X / (X - P_{ri} * P_{ri}), \quad P_{ri} < N_{ri} \text{ のとき} \quad (\text{式 6})$$

$$= N_{ri} + 1.0, \quad P_{ri} \geq N_{ri} \text{ のとき}$$

【0037】

式 6 は、図 2 に示す 1 つの計算機（CPU 1 台）における平衡状態の平均値に関して成立する関係から、以下のようにして導かれる。

【0038】

トランザクションは指数分布に従う時間間隔で到着する（ポアソン到着）とする。また、ディスク装置では待ちは生じない（装置が無限に存在する）ものとする。前述の処理中業務プロセス数、その内の CPU 系に存在する業務処理プロセス数、CPU 使用率もここでは平均値とし、これらを含めてすべての変数は計算機番号、推定状態を示す添字を省いて示す（例えば、処理中業務プロセス数は N 、CPU 系に存在する業務処理プロセス数は P 、CPU 使用率は R でそれぞれ示す）。また、以下で定義する 4 種の変数はトランザクション当たりの平均時間とする。

F : 処理時間

t : 純処理時間

s : CPU 使用時間

d : 入出力時間 ($t = s + d$)

更に、対象計算機からの要求で実行中の入出力数の平均を D とする ($N = P + D$)。

従来技術 4 の文献 228 頁の 8.3 式から、平衡状態の平均値について、

$$F = s(P + 1) + d = sP + t \quad (\text{式 61})$$

となり、入出力で待ちがないので、

$$d / s = D / R \quad (\text{式 62})$$

となる。また、同文献 228 頁の式 8.3 と式 8.1 の対比からも知られるように、

$$P+1=1/(1-R) \quad \text{から} \quad R=P/(1+P) \quad \text{式(63)}$$

となる。

式63を式62に代入し、 $d=t-s$ 、 $N=P+D$ を適用すると、

$$s=Pt/(P+D+PD)=Pt/(N(P+1)-P \cdot P) \quad \text{式(64)}$$

となり、式64を式61に代入すると、

$$\begin{aligned} F &= P \cdot Pt / (N(P+1) - P \cdot P) + t \\ &= N(P+1)t / (N(P+1) - P \cdot P) \\ &= Xt / (X - P \cdot P) \end{aligned}$$

となる。ここで、 $X=N(P+1)$ である。

したがって、伸長率Eは、

$$E=F/t=X/(X-P \cdot P)$$

となり、式6が導かれる。

また、伸長率Eは、NとCPU使用率Rを用いて次のように表すこともできる

$$E=N(1-R)/(N(1-R)-R \cdot R) \quad \text{式(65)}$$

式65は、式6にRとPの関係を表す式63を適用して得ることもできるし、従来技術4の文献の式8.1(Rを用いて処理時間を表現)から出発して、式62、式63を適用して得ることもできる(導出の記述は省略する)。しかし、Rについては式5で行ったのに相当する補正の手段がなさそうなので、本実施形態では補正の前にPに変換してから補正を受けるようにしてしまい、CPU使用率を測定した場合にも最終的な伸長率の式としてはNとPを用いる式6を使用するようにした。

なお、本実施形態で推定伸長率を求めるために用いる計算式(式6)は、このようにシステムの統計的平衡状態に関して成立するものであり、平衡状態がある程度の時間続くときに、その間のP、Nの平均値を知れば推定可能になるものである。現状を表す平衡平均値としては、過去の履歴に基づいて式3を用いた N_{ei} 、 P_{ei} が適当と考えられるが、Nについては正確な現在値 N_{pi} が知られているので、動的負荷配分の立場からはこれも反映すべく、前述した式4においては、この方針によりNの補正值を得ている。

【0039】

再び図4 (A) を参照して実行計算機選択手段7の残りの動作を説明する。

【0040】

採用する負荷指標によっては、現状における推定伸長率 E_{pi} の他に到着メッセージを計算機 i にスケジュールした場合の予測伸長率 E_{ni} が必要になる。あるいは E_{ni} だけを必要とすることもある。 E_{ni} も E_{pi} と同様に図4 (A) のステップ73で計算される。

【0041】

E_{ni} が必要であって到着メッセージ処理のジョブ特性を利用しない場合は、スケジュール後のトランザクション数 N_{ni} を $N_{ri} + w$ とし、スケジュール後のCPU系滞在プロセス数 P_{ni} を $P_{ri} + w$ として、式6と同様に E_{ni} を計算しテーブルの列T6に格納する。

【0042】

到着メッセージ処理のジョブ特性が推定可能でこれを利用する場合は、メッセージの種類などからその純処理時間 (CPU、ファイル装置という資源を実際に使用する時間の合計、言い換えると資源競合が全くない場合の処理時間) に占めるCPU時間の割合 C を推定し、これを用いて次の計算により、まず前記 P_{ni} を推定する。

【0043】

C_1 をスケジュール前における計算機 i 上における C の推定値とすると、 C_1 は式64における s/t に相当するので、式64に N 、 P の補正值を当てはめて、

$$C_1 = P_{ri} / (N_{ri} * (1 + P_{ri}) - P_{ri} \cdot P_{ri})$$

となる。 C_2 を計算機 i に、 $C = C_0$ である到着メッセージをスケジュールした場合の新ジョブミックスにおける C の推定値とする。平均が C_1 のジョブが N_{ri} 個存在し、そこへ C_0 のものが1個加わり総数は N_{ni} となるので、その平均値は、

$$C_2 = (N_{ri} * C_1 + C_0) / N_{ni}$$

となる。

式 6 4 から、 $s/t = P / (NP + N - P \cdot P)$ なので、この式をスケジュール後の状態に適用すると、 $s/t = C_2$ なので、 P_{ni} を y とおくと、

$$C_2 (N_{ni} \cdot y + N_{ni} - y^2) = y$$

となり、整理すると、次の 2 次方程式が得られる。

$$C_2 y^2 + (1 - C_2 \cdot N_{ni}) y - C_2 \cdot N_{ni} = 0 \quad (\text{式 7})$$

式 7 を y について解くことによって、スケジュール後の P_{ni} の推定値が得られる。定数項が負の値なので正の解と負の解が得られる。正の解を P_{ni} として採用する。そして、式 6 と同様にして E_{ni} を計算する（これを、以下では E_{ki} とする）。

【 0 0 4 4 】

ここで処理時間の伸長率とは、業務処理プロセスの応答時間、すなわち待ち時間も含む処理時間の、純処理時間に対する倍率を表す。伸長率 E_i は、計算機 i における業務処理プロセスの伸長率である。処理速度が同じ計算機ならば、同一の処理は伸長率の小さい計算機で実行した方が処理時間は短く、したがって応答時間を短くできることになる。推定伸長率は、当該計算機上で実行中のプロセスの集まり（ジョブミックス）の、動作中の群としてのプログラム特性（CPU 使用特性だけでなく、CPU-I/O 使用特性を含む）を反映している。しかも、式 6 を用いると、実行中の個々のジョブの特性を知る必要がなく、動作中に観測可能なデータだけから得ることが可能なところに特徴がある。基本的に、従来技術 4 の考え方の系列に属し、式 6 は式 2 の拡張・変形により得られるが、当方式は現時点のシステム状況（ジョブミックス特性）を反映可能にし、かつ、CPU 系での滞在時間だけでなく入出力も含めた全処理時間（応答時間）を対象にして、精度・ダイナミック性を向上させている。ただし、式 6 は平衡平均値に関する理論に基づいているので、短期的な状況の把握法としては 100 パーセントの信頼性があるとは言えない。

【 0 0 4 5 】

負荷指標の値の計算（図 4 A の 7 4）では、負荷データ $T_2 \sim T_6$ を用いて各計算機について負荷指標の値を計算する。負荷指標としては図 5 に示すように 8 種類（名称として L で始まる）が考えられ、実施システムではこの内の一種類を

選べばよい。図中に示した式による計算で結果を得てテーブルの推定負荷の列 T 7 に格納する。いずれも、小さい値を持つ計算機ほどスケジュール先として望ましいことになる。どの時点の負荷を考えるかについて、到着メッセージのスケジュール前（この負荷を L_p と表記する）／後があり、さらに、スケジュール後の場合に到着メッセージのジョブ特性を未知とする（負荷を L_a と表記）か、推定可能とする（負荷を L_k と表記）かがありうる。これら 3 ケース各々について、伸長率そのものを負荷指標と捉える（ L_{x1} と表記）こともでき、推定伸長率に処理中業務処理プロセス数 N_{ri} または N_{ni} を乗じたものを負荷指標とする（ L_{x2} と表記）こともできる（ x は p 、 a または k である）。後者は計算機上の個々のトランザクションの推定伸長率の総和という性格をもつ。さらに、スケジュールによる負荷の増加という観点から、上記の総和のスケジュール前後における増分を負荷指標とする（ L_{x3} と表記）こともできる。伸長率の増分最小という選択は、システム全体にとって当スケジュールによる応答時間総和の増加を最小にする選択になり、結果として平均応答時間を最小化できると期待できる。到着メッセージのジョブ特性が相当の精度で推定可能な場合は、理論通り、 L_{k3} を採用するのが最も良い結果を期待できる。ジョブ特性推定の精度が期待できない場合は L_{k3} の選択は危険であり、平均応答時間最小という点からは L_{a2} を採用するのがよい。

【0046】

推定値の補正（図 4 A の 7 2）、推定伸長率の計算（7 3）、負荷指標の値の計算（7 4）は、入力メッセージを処理可能なすべての計算機に関して行い、推定負荷を得てテーブルの列 T 7 に格納しておく。

【0047】

実行すべき計算機の決定（図 4 A の 7 5）では、テーブルの列 T 7 に格納されている各計算機の前記推定負荷をサーチし、推定負荷が最小の計算機（計算機 j とする）を選択する。次に、メッセージの送付（図 4 A の 7 6）では、選択された計算機 j に対して入力メッセージを送付して処理開始を促すように、通信手段 21 に指令する。

【0048】

次に、本実施の形態の効果について説明する。

【 0 0 4 9 】

本実施の形態では、中継配分装置 2 の上で全計算機の負荷データをリアルタイムで管理し、またすべての処理要求メッセージを直接受け取り、直ちに、前記負荷データに基づいて各計算機における伸長率を計算し、その時点で最適な負荷指標値をもつ計算機に処理要求メッセージの処理を依頼するように構成されているため、集中的な制御が実現でき、オーバヘッドの少ない、かつ、良質な負荷配分を実現することができる。

【 0 0 5 0 】

次に、本発明の第 2 の実施の形態について図面を参照して詳細に説明する。

【 0 0 5 1 】

図 6 を参照すると、本発明の第 2 の実施の形態は、第 1 の実施の形態に対し、構成として、中継配分装置 2 をもたず端末群 5 は通信網 5 1 を介して直接に計算機群 1 に接続されている点と、計算機間を接続する交換・蓄積機構 1 0 が追加されている点が異なる。これに伴い、計算機 1 i は、負荷データ測定手段 A 1 i 1 と、トランザクション処理手段 1 i 2 と、通信手段 1 i 3 に加えて、負荷データ記憶手段 6 i と、実行計算機選択手段 7 i と、負荷データ測定手段 B 8 i とを含む。これらの各手段を他の手段と共にソフトウェア的に実現する場合、第 1 の実施の形態と同様にその実現用プログラムが図示しない記録媒体に記録されて提供される。

【 0 0 5 2 】

1 つのトランザクションの処理の概略は次のようになる。トランザクション処理要求であるメッセージは端末群 5 に属する端末装置から送り先を指定して送り出され、通信網 5 1 を経由して指定された計算機で受け取られる。受け取った計算機はそのメッセージを自分で処理するか他に依頼するか、依頼するとしたらどの計算機にするかを決定し、依頼する場合は交換・蓄積機構 1 0 を介して依頼先計算機に該メッセージを送る。処理を行う計算機はトランザクション処理を実行し、応答メッセージを作成して要求元端末に返す。

【 0 0 5 3 】

ここで、上記の手段はそれぞれ概略つぎのように動作する。

【0 0 5 4】

計算機 1 i 上の負荷データ測定手段 A 1 i 1 は、一定時間ごとに自身の属する計算機 1 i の負荷データを測定し、負荷データ測定手段 B 8 i でこれを加工して推定データとして負荷データ記憶手段 6 i に格納する、と共に交換・蓄積機構 1 0 により他のすべての計算機に通知する。また、同様に各計算機のトランザクション処理開始および終了を相互に通知し合う。これらによって、各計算機上の負荷データ記憶手段 6 x には全計算機の最新の負荷データが保持される。端末から来た処理要求メッセージは実行計算機選択手段 7 i が受け、負荷データ記憶手段 6 i に記憶されている各計算機の推定負荷データから各計算機における推定伸長率を求め、これに基づいて実行すべき計算機を決定し、自身で実行する場合にはトランザクション処理手段 1 i 2 に渡し、他の計算機に実行させる場合は該計算機のトランザクション処理手段 1 j 2 に向けて、交換・蓄積機構 1 0 を経由して処理要求メッセージを送付する。トランザクション処理手段 1 x 2 は、処理要求メッセージを受け取ると、自身の管理下の業務処理プロセスを割り当てて処理を行わせる。業務処理プロセスは、プログラムの実行のために CPU の使用とファイル装置 4 上のファイルへのアクセスを繰り返し、処理が終了すると応答メッセージを作成し、通信手段 1 x 3 を介して要求元端末に送る。

【0 0 5 5】

次に、図 6 及び図 4 B のフローチャートを参照して本実施の形態の全体の動作について詳細に説明する。

【0 0 5 6】

トランザクション処理手段 1 i 2 の動作は、処理要求メッセージを受け取るのが自分または他の計算機上の実行計算機選択手段 7 x からである点を除くと第 1 の実施の形態と同一である。負荷データとして管理するデータも第 1 の実施の形態と同一で図 3 に示すものであるが、各計算機上に負荷データ記憶手段 6 x として、全計算機に関する同一内容のものを保持する。計算機 1 i 上の負荷データ測定手段 A 1 i 1 は、一定時間ごとに自身の属する計算機 1 i の負荷データとして、CPU 系に滞在する業務処理プロセス数 P i あるいは CPU 使用率 R i、およ

びその時点での処理中業務処理プロセス数 N_i を測定する。そして、負荷データ測定手段 $B8_i$ がこれを加工して、推定データとして、 P_{ei} あるいは R_{ei} 、および N_{ei} を求め、前記 R 方式なら R_{ei} から P_{ei} を計算し、列 $T3$ に N_{ei} を列 $T4$ に P_{ei} をそれぞれ i 番目の値として格納する。同時に N_{ei} 、 P_{ei} の値を他のすべての計算機に交換・蓄積機構 10 を介して送り、負荷データ記憶手段 $6x$ の内容を更新させる。推定データの計算方法は第 1 の実施の形態におけるのと同じである。また、計算機 $1i$ 上で処理中のトランザクション現在数 N_{pi} (テーブル上の列 $T2$) に関しては、負荷データ測定手段 $B8_i$ が、計算機 $1i$ でのトランザクション処理開始およびトランザクション処理終了の通知をトランザクション処理手段 $1i2$ から受けて更新・保持する、と共に他のすべての計算機に送り負荷データ記憶手段 $6x$ の内容を更新させる。

【0057】

図 $4B$ は、本発明の第 2 の実施の形態の実行計算機選択手段 7 の動作を示すフローチャートである。端末から処理要求メッセージが到着した計算機 $1i$ において実行計算機選択手段 $7i$ が実行される。到着メッセージを受けた通信手段 $1i3$ に起動されて該メッセージを渡され (図 $4B$ の 71)、負荷データ記憶手段 $6i$ に記憶されている各計算機の負荷データから各計算機における推定伸長率を求めて (図 $4B$ の 72 , 73) 負荷指標の値を計算し (図 $4B$ の 74)、これに基づいて自計算機 ($1i$) で実行すべきかどうか判断し (図 $4B$ の 751)、自計算機ですべきでないなら実行する計算機を決定し (図 $4B$ の 75)、選択された計算機に向けて、処理要求メッセージを交換・蓄積機構 10 を介して送付させる (図 $4B$ の 76)。自計算機で実行すべきなら自分を選択し (図 $4B$ の 752)、処理を指示する (76)。実行計算機選択手段 $7x$ の動作として第 1 の実施の形態と論理的に異なるのは、自計算機で実行するか否かの判断のところだけである。推定値の補正 (図 $4B$ の 72)、推定伸長率の計算 (73)、負荷指標の値の計算 (74) は、入力メッセージを処理可能なすべての計算機に関して行い、推定負荷を得てテーブルの列 $T7$ に格納しておく。自計算機で実行するか否かの判断 (751) には、まず、自計算機の現推定伸長率 E_{pi} を用い、これが閾値 (小さめに 1.3 程度がよい) 以下であったならば、自計算機で実行することに

する。そうでない場合、負荷データ記憶手段 6 i 上の推定負荷 T 7 に基づいて、自計算機の推定負荷が最小でなくても、最小負荷の計算機との差が小さければ自計算機で実行するようにする。較差の大小の判断は閾値による。負荷指標として伸長率総和の増分以外の 6 種類のいずれかを採用する場合は、較差の判断は倍率閾値によるのがよく、1. 3 倍から 1. 5 倍程度がよいようである。すなわち、自計算機の E_{pi} が E_p が最小である計算機 j の E_{pj} の 1. 3 倍以内であったら自計算機で実行する、などである。負荷指標として伸長率総和の増分を採用する場合は、2 台の計算機の負荷指標間の差をシステム内の全トランザクション数で除したものが閾値を越えるか否かで判断するのがよい。すなわち、実行中の全トランザクションの平均伸長率の増分の程度によって判断する。この場合の閾値は 0. 0 2 程度がよい。ここで、閾値の値は、負荷の測定間隔が長い場合には大きくした方がよい。これは、測定間隔が長い場合は負荷指標の推定値の信頼性が低くなるので、トランザクション転送を行う頻度が少なくなるような安全サイドの選択をした方がいいからである。到着計算機で処理を実行してしまうことを優先するのは、他の計算機に転送するには転送元・転送先の双方にオーバーヘッドがかかり、また、対象トランザクション自身の処理時間に遅延をもたらすからである。閾値を用いた判断を入れないと、ほとんどすべてのメッセージのトランザクション処理を他の計算機へ依頼する結果になる可能性が高い。

【 0 0 5 8 】

自計算機で処理することに決定したらトランザクション処理手段 1 i 2 にメッセージを引き渡し処理を依頼する。到着計算機で処理すべきでないとなったときは、実行すべき計算機の決定（図 4 B の 7 5）で、テーブルの列 T 7 に格納されている各計算機の推定負荷をサーチし、推定負荷が最小の計算機（計算機 j とする）を選択する。そして、メッセージの送付（図 4 B の 7 6）で、選択された計算機 j のトランザクション処理手段 1 j 2 に対して、入力メッセージを交換・蓄積機構 1 0 を介して送付し、処理開始を促す。

【 0 0 5 9 】

以上において、負荷データ記憶手段 6 は同一内容のものが各計算機上に保持されるとし、各計算機で自身に関して測定／計算後に他のすべての計算機に交換・

蓄積機構 10 を介して通知するとしていたが、交換・蓄積機構 10 がある程度の容量をもち主記憶程度に速い蓄積機構を備えるなら、前記負荷データ記憶手段 6 の一部は、システム共用のものとして交換・蓄積機構 10 の上に格納し保持することもできる。負荷データ更新のオーバヘッドの観点から、この構成の方が望ましい。この場合、テーブルの列 T1～T4 は交換・蓄積機構 10 に保持し、各計算機はメッセージが到着した際に、ここから引き出したデータに基づいて推定伸長率、推定負荷などを計算し、処理を実行すべき計算機を決定すればよい。また、交換・蓄積機構 10 が前述の条件を満たす場合、他の計算機に処理を依頼することになったときには、メッセージそのものを直接送付するのでなく、メッセージは交換・蓄積機構に格納し、依頼の通知だけを相手に送るように構成することもできる。この場合、受け取り側の計算機は、都合の良いときに非同期的に交換・蓄積機構から取り出すことになる。

【0060】

次に、本実施の形態の効果について説明する。

【0061】

本実施の形態では、特別な中継配分装置を備えなくてもよい。システム全体として低コストで構成することができる。集中制御による負荷分散はできないが、処理要求メッセージは端末からの指定により送付された先の計算機で、その計算機及び他の計算機の負荷状況データに基づいて伸長率を計算し、転送のオーバヘッドも考慮した上で、その時点で最適な実行計算機を決定し、その計算機上で実行させるように構成されているため、集中制御である第 1 の実施の形態よりは落ちるが、分散制御下としては高い応答性能を実現できる。

【0062】

図 8 及び図 9 に示すグラフは、本実施の形態におけるような分散制御下における負荷分散の効果を、シミュレーション評価によって確認した結果である。トランザクションとしては、純処理時間（450 ミリ秒）に占める CPU 時間の割合が平均 5 % のものと平均 60 % のものの 2 種類が、7 対 3 の割合で到着するとした。計算機は 8 台あり、各計算機への到着はランダムで、平均としては等しい到着率になるように設定した。横軸は到着率に比例する負荷率を示し、縦軸は図 8

では得られた平均応答時間（ミリ秒）であり、図 9 では応答時間のばらつき（標準偏差）である。それぞれのグラフ曲線は負荷分散方式に対応しており、実線のものが本実施の形態に関係する。NC 方式は負荷分散をせず、到着したものをそのまま処理する。MPL 方式は、処理中トランザクション現在数を負荷指標とする動的制御で、これが到着計算機より 2 以上小さい計算機が存在したら、最小の計算機に転送し処理させる。La 2、Lk 3 はそれぞれ本実施の形態における推定伸長率に基づく負荷指標を用いた動的制御に対応する。これらの結果から、平均応答時間について、静的確率的配分としては最適であるはずの NC 方式よりも動的制御は大幅によいことが分かり、特に推定伸長率に基づく方式は従来多く用いられている実行中トランザクション数に基づく方式よりも優れていることが示され、また、応答時間のばらつきについても同様な傾向が、より顕著に現れていることが分かる。このような差は負荷率が高いときに、より顕著である。

【0063】

次に、本発明の第 3 の実施の形態について図面を参照して詳細に説明する。

【0064】

図 7 を参照すると、本発明の第 3 の実施の形態は、第 2 の実施の形態に対し、構成として、中継仮配分装置 25 をもち、端末群 5 は中継仮配分装置 25 を経由して計算機 11 ～ 1n に接続されている点だけが異なる。

【0065】

1 つのトランザクションの処理の概略は次のようになる。トランザクション処理要求であるメッセージは端末群 5 に属する端末装置から送り出され、前記中継仮配分装置 25 に渡される。中継仮配分装置 25 は受け取ったメッセージを処理する計算機 1i を仮決定し、該計算機に該メッセージを送る。受け取った計算機 1i はそのメッセージを自分で処理するか他に依頼するか、依頼するとしたらどの計算機にするかを決定し、依頼する場合は交換・蓄積機構 10 を介して依頼先計算機に該メッセージを送る。処理する計算機はトランザクション処理を実行し、応答メッセージを作成して要求元端末に返す。ここで、第 2 の実施の形態に対して追加された中継仮配分装置 25 は、端末からのメッセージを受けて仮配分先計算機を決定して送付するが、基本的に、配分は詳細な動的情報に基づかない、

静的／準静的な手法によって行われる。

【 0 0 6 6 】

次に、本実施の形態の全体の動作について詳細に説明する。

【 0 0 6 7 】

中継仮配分装置 2 5 は、第 1 の実施の形態における前記中継配分装置 2 と同様に、端末装置群 5 から送り出されるすべての処理要求メッセージを受け取り、これを渡すべき計算機を決定して送付する。中継仮配分装置 2 5 における静的／準静的な仮配分方式として次の 3 種類が想定される。実施に当たっては、このうちいずれか 1 種類か、あるいはこれらを組み合わせた方式を選択する。これら以外であっても、計算機からの負荷データの収集が少なく、実行のオーバーヘッドも小さい配分方式なら採用可能である。

- (1) 端末のグループ分けによる配分
- (2) 到着順に、巡回的に各計算機へ配分
- (3) 実績データに基づく確率的配分

【 0 0 6 8 】

(1) 端末グループ分けによる配分では、メッセージ発生元の端末によって配分先の計算機を固定的に定めておく。すなわち、端末群を計算機 1 で処理するグループ、計算機 2 で処理するグループ、のように予めグループ分けしておき、どの端末から来たかによって行く先を機械的に決定する。第 2 の実施の形態とほとんど同じ方式になるが、本方式では端末群と計算機との対応関係を中継仮配分装置で集中的に管理できるので、過去の実績に応じて、例えばシステム立ち上げの度ごとに、長期的には負荷バランスのとれるグループ分けに設定し直すなどを容易にできる。

【 0 0 6 9 】

(2) 到着順に巡回的に各計算機へ配分では、中継仮配分装置に到着した最初のメッセージは計算機 1 へ、次は計算機 2 へ、と順次配分し、最後の計算機 n に配分した次のメッセージは再び計算機 1 へ、と巡回的に配分する。特に大部分のメッセージ処理のジョブ特性が同一クラスに属するような場合、短期的にも負荷をバランスさせる効果が期待できる。

【0070】

(3) 実績データに基づく確率的配分では、各計算機に配分するメッセージ数の比率を計算機ごとに設定し、短期的にもこの比率を守るように配分をする。各計算機から負荷状況のデータを1秒ごと、10秒ごとなどに定期的に受け取り、負荷がアンバランスであったなら、バランスさせるように個々の計算機への配分比率を上下させ、以後はこの配分比率に基づいて配分を行うようにする。

【0071】

図7におけるトランザクション処理手段1x2の動作、各計算機上にある負荷データ記憶手段6xの内容は第2の実施の形態と同一である。計算機1x上の負荷データ測定手段A1x1、負荷データ測定手段B8xも第2の実施の形態におけるのと同じ動作をするが、それに加えて負荷データ測定手段Bは、中継仮配分装置25が前述の配分方式(3)を採用する場合、1秒、10秒などの間隔で負荷データの概要を中継仮配分装置25に送る。本発明の第3の実施の形態の実行計算機選択手段7の動作は、第2の実施の形態におけるのと同じであり、図4Bのフローチャートで示される。

【0072】

以上において、負荷データ記憶手段6は同一内容のものが各計算機上に保持され、各計算機で自身に関して測定/計算後に他のすべての計算機に交換・蓄積機構10を介して通知するとしていたが、負荷データ記憶手段6はシステム共用のものとして交換・蓄積機構10の上に格納し保持することもできる。負荷データ更新のオーバーヘッドの観点から、この構成の方が望ましい。また、他の計算機に処理を依頼することになったときには、処理対象メッセージそのものを直接送付するのではなく、メッセージは交換・蓄積機構に格納し、依頼の通知だけを相手に送るように構成することもできる。以上の点に関しても、第2の実施の形態におけるのと同様である。

【0073】

次に、本実施の形態の効果について説明する。

【0074】

本実施の形態では、中継配分装置として、限定された機能だけをもち、計算機

群 1 1 ~ 1 n からの情報収集量・頻度も小さいものを備えるだけでよいので、比較的 low コストで全体システムを構成できる。機能の限定された中継仮配分装置であるが、ここで準静的にとはいえ負荷の適切な仮配分を行うように構成されているため、中継仮配分装置が存在しない場合と比較して、応答性能（平均、ばらつき共）を向上させることができ、また、計算機に到達してから行われる負荷バランスのための転送の頻度を大幅に減少させることができる。

【 0 0 7 5 】

図 8 及び図 9 は、第 2 の実施の形態の効果の説明で前述した条件の下で、シミュレーション評価により得られたものであり、第 3 の実施の形態の結果が点線のグラフとして含まれている。グラフ曲線はそれぞれ負荷分散方式に対応している。R__NC は、仮配分として前記（2）到着順に巡回的に配分を実施して、仮配分先の計算機でそのまま処理を実行させたものである。R__Lk 3 は、仮配分を同じく（2）で行い、仮配分先の計算機で前記 Lk 3 を負荷指標とする負荷配分を行った結果である。これから、準静的な配分だけでも静的な配分である NC 方式よりも応答性が大幅に向上することが分かり、さらに推定伸長率に基づく負荷分散を組み合わせることにより、第 2 の実施の形態によるよりも応答性を向上させることが理解できる。

【 0 0 7 6 】

【発明の効果】

本発明によれば、複数の計算機が負荷分担してトランザクション処理を実行するシステムにおいて、短期的にも計算機間の負荷をバランスさせ、全体として応答時間の平均及びばらつきを小さく保つことができる。その理由は、計算機の処理時間の伸長率をベースとする新規な負荷指標に基づいて処理要求を動的に配分しているためである。つまり、基本的な性能指標として処理時間の伸長率は、当該計算機上で実行中のプロセスの集まり（ジョブミックス）の、動作中の群としてのプログラム特性（CPU 特性だけでなく、CPU-I/O 使用特性を含む）を反映しており、対話型処理にとっては最適で、かつ、適用性が広いためである。

【 0 0 7 7 】

また本発明によれば、各計算機における処理時間の伸長率を、一定時間ごとの各計算機の処理中トランザクション数、CPU系に滞在する業務処理プロセス数、CPU使用率といった、動作中に観測可能なデータから導き出せるようにしたため、個々の処理要求のジョブ特性に関する先験的知識なしに、低オーバーヘッドで実測可能な負荷データだけに基づいて、伸長率をベースとした負荷指標に基づく負荷分散が実現できる。

【0078】

また一定時間ごとに測定した負荷データの系列を総合的に用いて各計算機の負荷状況を推定し、あるいは各計算機におけるトランザクション処理の開始・終了に応じて各計算機の処理中トランザクション現在数を常に把握しておき、この処理中トランザクション現在数を用いて推定負荷状況データを補正する構成にあっては、推定負荷状況の推定精度が高まり、ひいては伸長率の推定精度、負荷分散の精度をより向上させることができる。

【0079】

推定伸長率をベースとした各計算機の負荷指標として、推定伸長率そのものを負荷指標としたり、当該計算機へ新規にトランザクションを割当て前あるいは割当て後における総推定伸長率（当該計算機における前記処理時間の推定伸長率に当該計算機の処理中トランザクション数を乗じた値）を用いたり、また、当該計算機へ新規にトランザクションを割り当てた後における前記総推定伸長率と、割当て前における前記総推定伸長率との差を用いる構成にあっては、個々の計算機のあるいはシステム全体としての負荷の程度を表現する負荷の指標に基づいた負荷配分が可能となり、個々のトランザクション自体の処理時間の最短化だけでなく、その割当てが他に及ぼす影響まで考慮した、システム全体としての最適化も可能となる。

【0080】

本発明は、集中的に動的に個々の処理要求を配分する構成のシステムにも、静的／準静的に配分をされてしまった後で受けた計算機が配分の修正という位置づけで処理の転送を行うことになるという構成のシステムにも適用が可能である。シミュレーションを用いた性能評価結果が、前述のように図8、図9に示されて

いる。

【図面の簡単な説明】

【図 1】

本発明の第 1 の実施の形態の構成を示すブロック図である。

【図 2】

任意の計算機 i から見たシステムのモデルである。

【図 3】

負荷指標の値を計算するために用いるデータの一覧表である。

【図 4】

(A) 第 1 の実施の形態における実行計算機選択手段の動作を示す流れ図である。

(B) 第 2 あるいは第 3 の実施の形態における実行計算機選択手段の動作を示す流れ図である。

【図 5】

本発明で用いる 8 種類の負荷の指標を示す一覧表である。

【図 6】

本発明の第 2 の実施の形態の構成を示すブロック図である。

【図 7】

本発明の第 3 の実施の形態の構成を示すブロック図である。

【図 8】

平均応答時間（本方式のシミュレーション結果）のグラフである。

【図 9】

応答時間の標準偏差（本方式のシミュレーション結果）のグラフである。

【符号の説明】

1 計算機群

1 0 交換・蓄積機構

1 1 ~ 1 n 計算機 1 ~ n

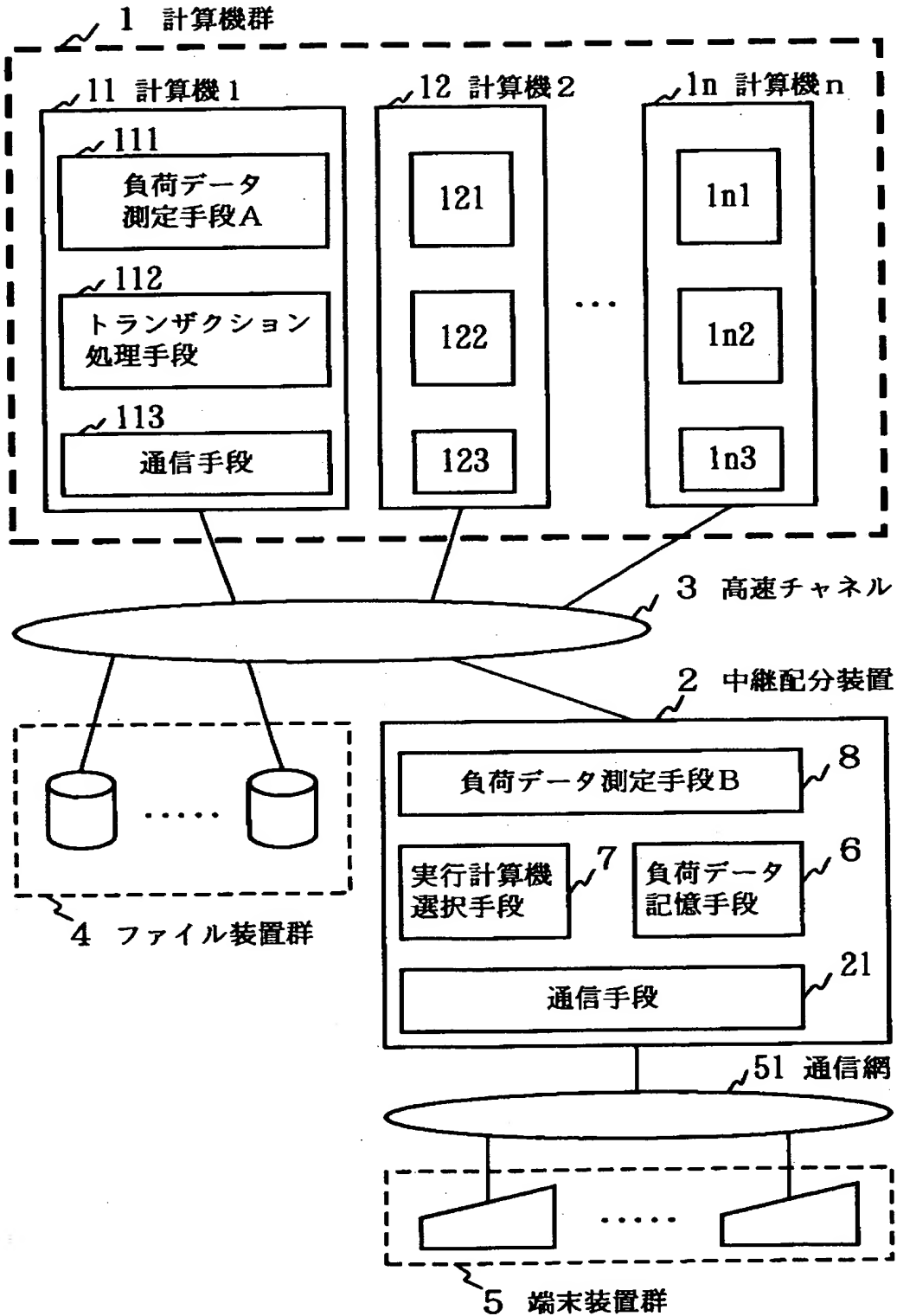
1 x 1 計算機 x 上の負荷データ測定手段 A

1 x 2 計算機 x 上のトランザクション処理手段

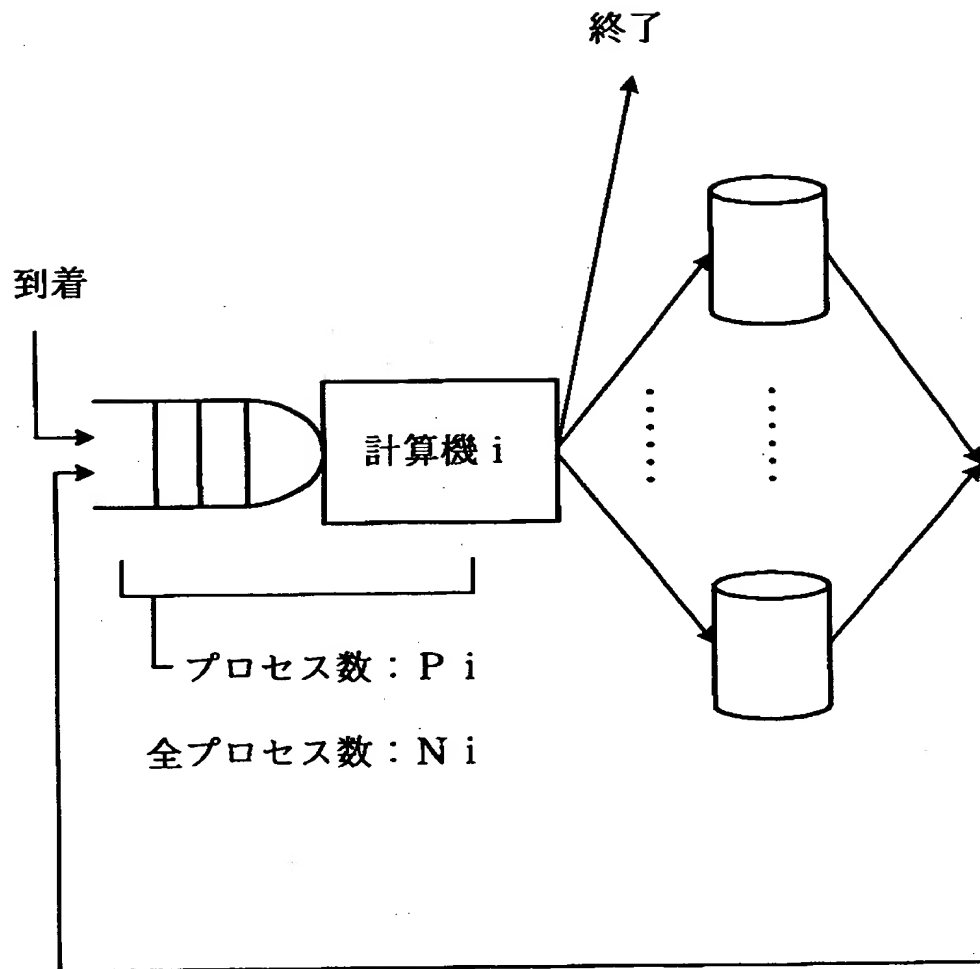
- 1 x 3 計算機 x 上の通信手段
- 2 中継配分装置
 - 2 1 通信手段
 - 2 5 中継仮配分装置
- 3 高速チャネル
- 4 ファイル装置群
- 5 端末装置群
 - 5 1 通信網
- 6 負荷データ記憶手段
 - 6 x 計算機 x 上の負荷データ記憶手段
- 7 実行計算機選択手段
 - 7 x 計算機 x 上の実行計算機選択手段
- 8 負荷データ測定手段 B
 - 8 x 計算機 x 上の負荷データ測定手段 B

【書類名】 図面

【図 1】



【図2】

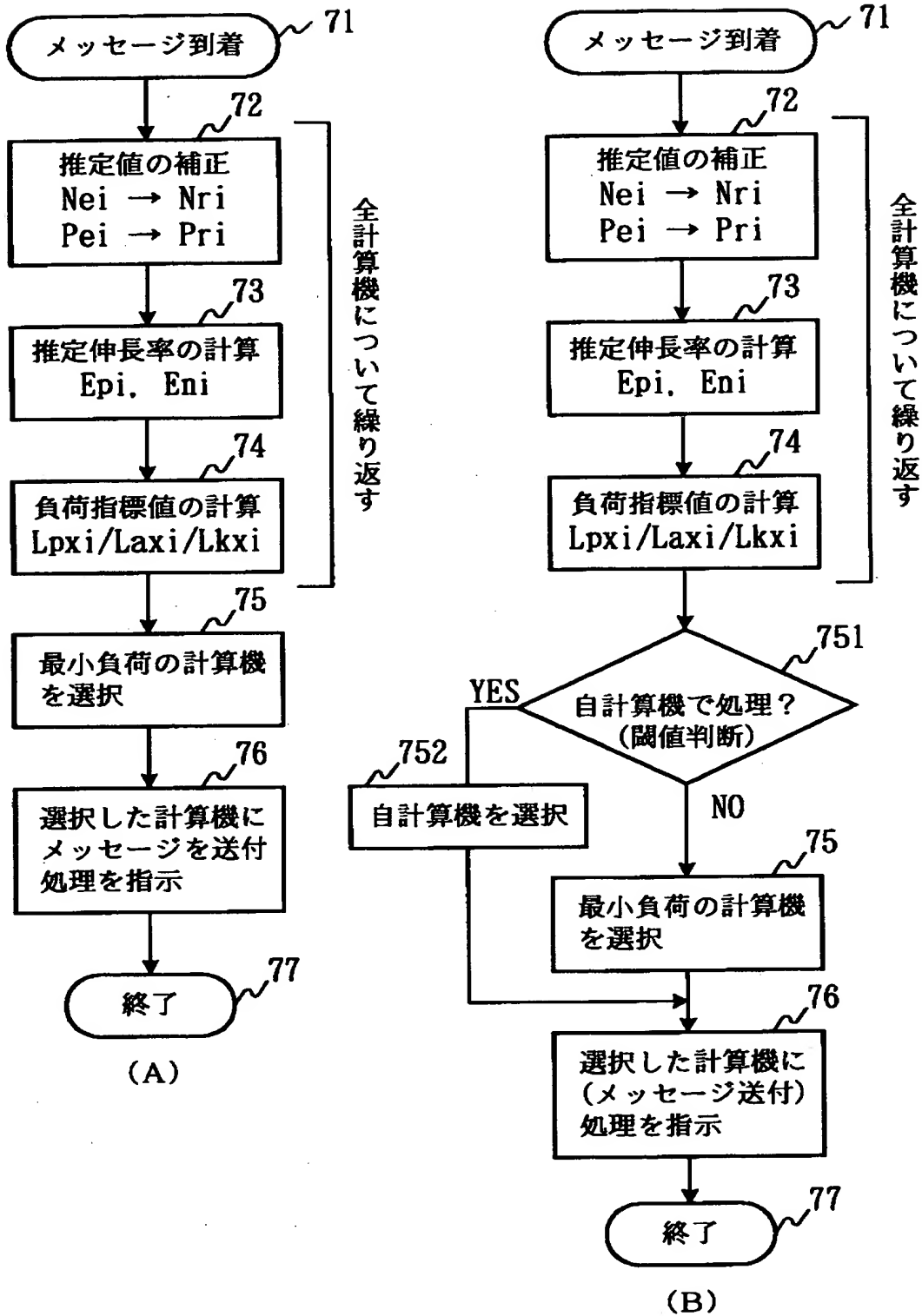


【図3】

T1	T2	T3	T4	T5	T6	T7
計算機 #	現 t r 数	推定 t r 数	推定 t r 数 / CPU 系	推定伸長率	配分後 推定伸長率	推定負荷
1	Np1	Ne1	Pe1	Ep1	En1	Lxy1
2	⋮	⋮	⋮	⋮	⋮	⋮
i	Npi	Nei	Pei	Epi	Eni	Lxyi
⋮	⋮	⋮	⋮	⋮	⋮	⋮
n	Npn	Nen	Pen	Epn	Enn	Lxyn

$$x = p/a/k \quad y = 1/2/3$$

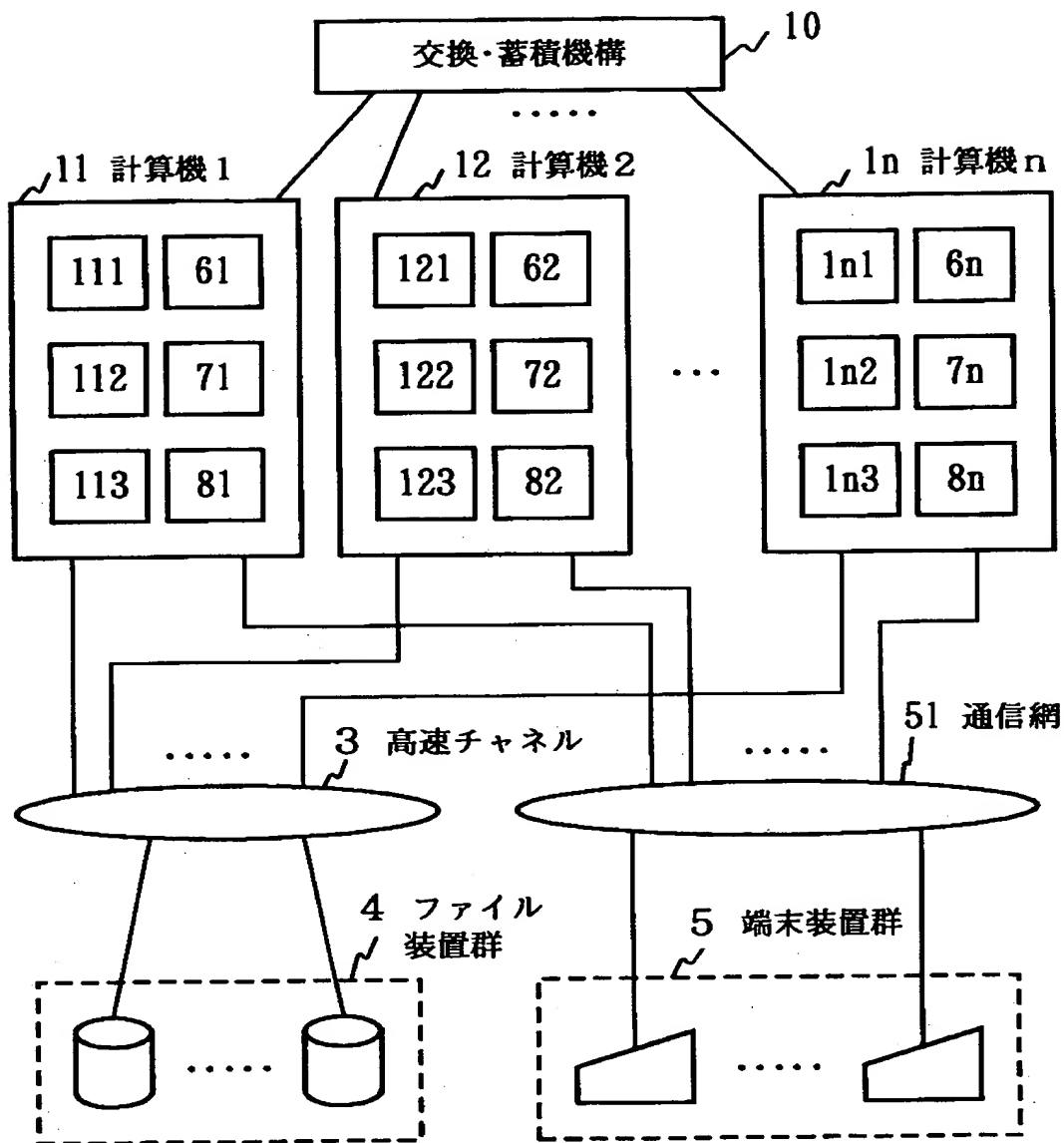
【図4】



【図5】

使用伸長率 負荷指標	スケジュール前 推定伸長率	スケジュール後推定伸長率	
		特性未知	特性推定
伸長率そのもの Tr 数 × 伸長率 総伸長率の増分	Lp1=Epi	La1=Eni	Lk1=Eki
	Lp2=Nri × Epi	La2=Nni × Eni	Lk2=Nni × Eki
	—	La3=La2 - Lp2	Lk3=Lk2 - Lp2

【図6】



61, 62, ..., 6n : 負荷データ記憶手段

71, 72, ..., 7n : 実行計算機選択手段

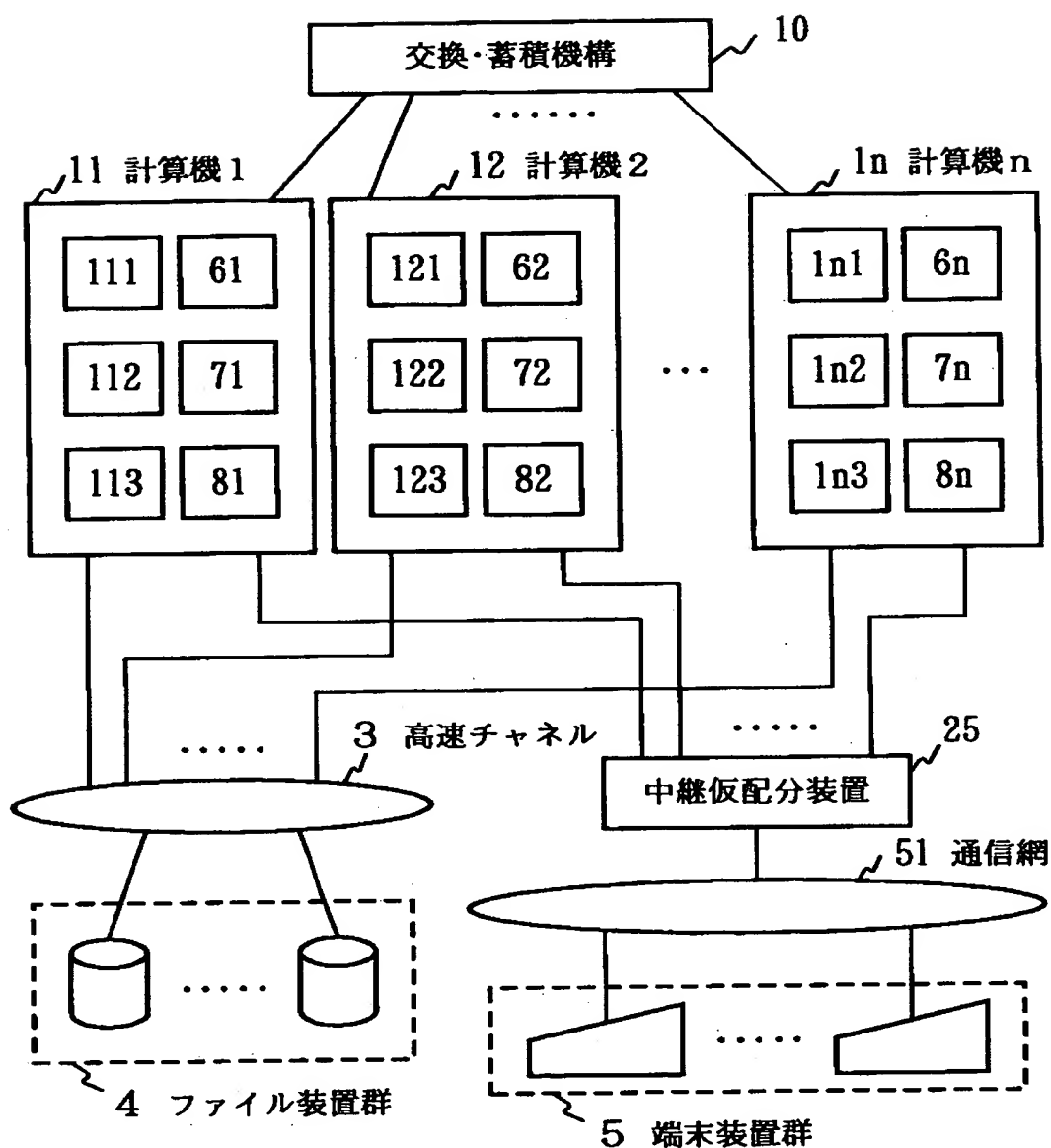
81, 82, ..., 8n : 負荷データ測定手段B

111, 121, ..., 1n1 : 負荷データ測定手段A

112, 122, ..., 1n2 : トランザクション処理手段

113, 123, ..., 1n3 : 通信手段

【図7】



61, 62, ..., 6n : 負荷データ記憶手段

71, 72, ..., 7n : 実行計算機選択手段

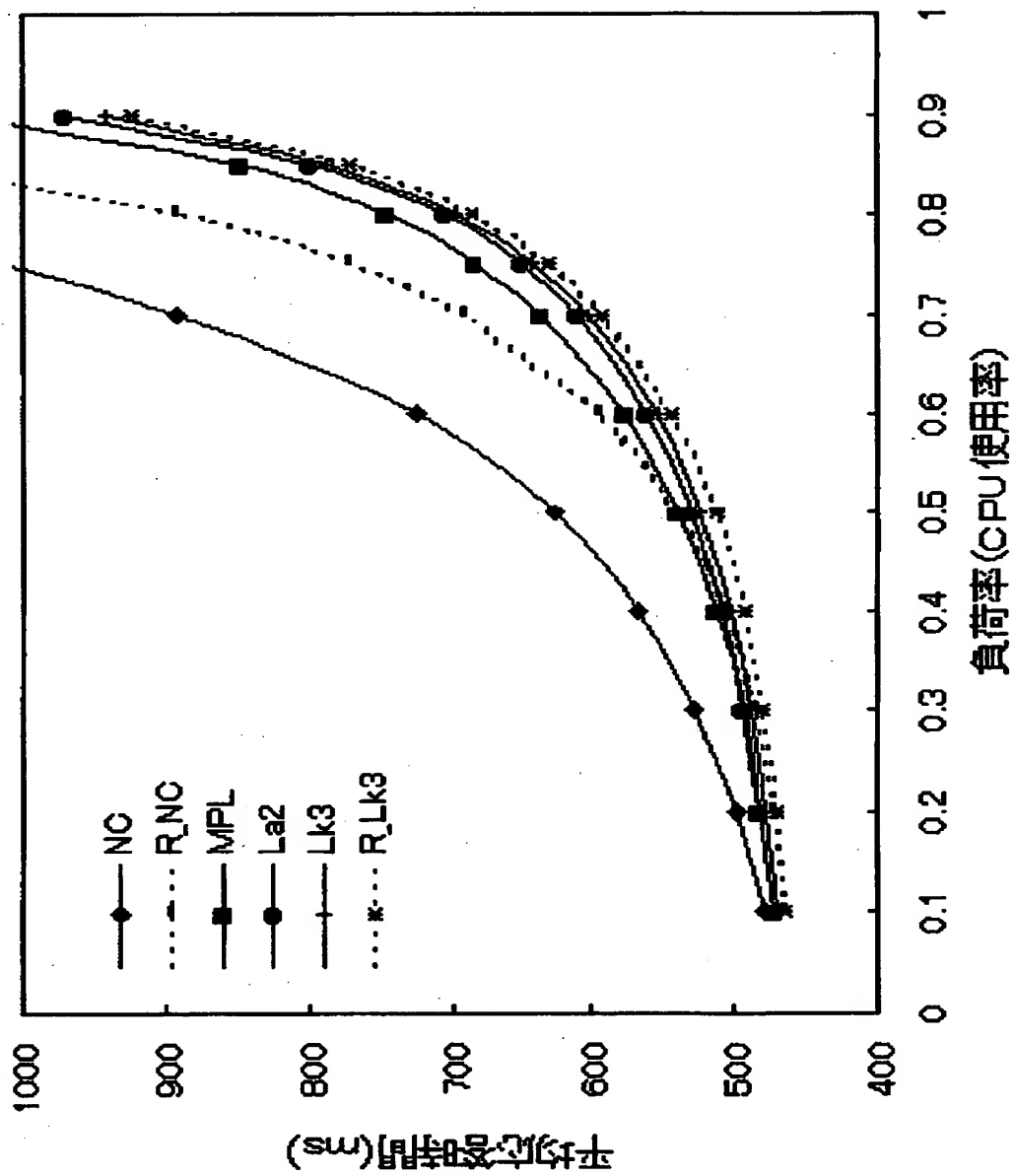
81, 82, ..., 8n : 負荷データ測定手段B

111, 121, ..., 1n1 : 負荷データ測定手段A

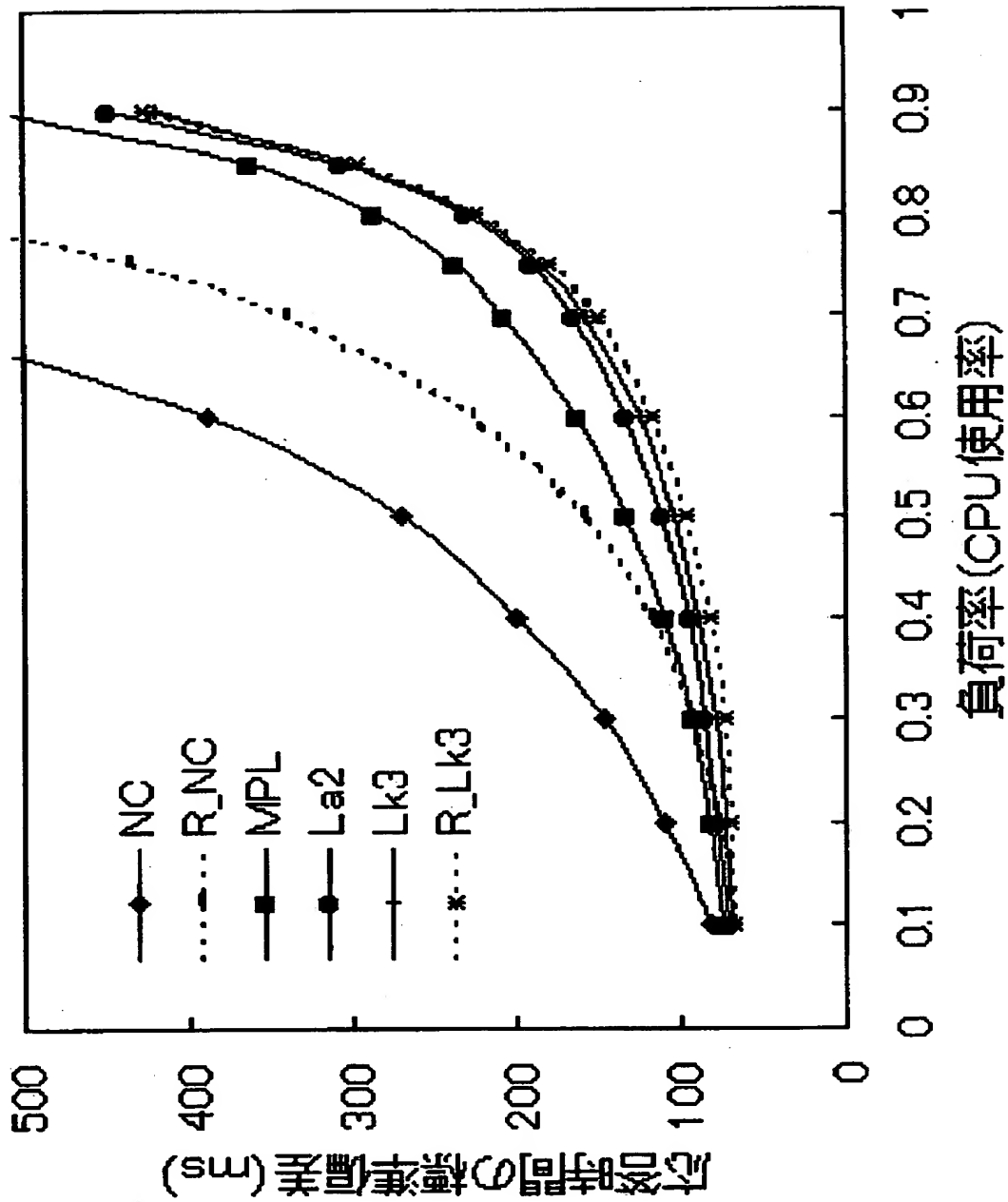
112, 122, ..., 1n2 : トランザクション処理手段

113, 123, ..., 1n3 : 通信手段

【図8】



【図9】



【書類名】 要約書

【要約】

【課題】 トランザクション処理要求を発生する端末装置群と要求の処理を負荷分担型により実行する複数の計算機からなるシステムにおいて、各計算機の負荷を動的にバランスさせることにより、全体の応答性能を良くする。すなわち、応答時間の平均とばらつきを小さくする。

【解決手段】 一定時間ごとに各計算機の処理中トランザクション数、およびCPU系に滞在するプロセス数またはCPU使用率を測定し、また、正確な処理中トランザクション現在数を常に把握し、これらの測定値を元に各計算機における処理時間の伸長率を推定し（73）、この推定伸長率に基づいて負荷指標の値を求め（74）、負荷の低い計算機へトランザクションの配分を行う（75、76）。各計算機上で実行中のジョブミックスの特性まで考慮した動的配分が実現でき、応答性能を上げられる。

【選択図】 図4

出 願 人 履 歴 情 報

識別番号 [000004237]

1. 変更年月日	1990年 8月29日
[変更理由]	新規登録
住 所	東京都港区芝五丁目7番1号
氏 名	日本電気株式会社